# Semantic Wiki for Collaborative Annotation of Chemical Biology Data

**Ajay[1], A. Zavaleta[2], A. Frazin[2], S. Chintala[2], and H. Cheung[2]**

[1]Cheminformatics Research Centers, National Human Genome Research Institute
[2]Division of Computational Bioscience, Center for Information Technology

**National Human Genome Research Institute**

**CIT** Center for Information Technology

## Abstract

As a central component of the Molecular Libraries Program (MLP), a NIH Roadmap project, the Molecular Libraries Program Center Network (MLPCN) is an effort to build high-throughput based, chemo-biological screening to test and understand the interaction between small molecular chemical structures and biological targets.

Resources for exchanging data and information have been ubiquitous in biology. The effort to enable cross-querying across databases with different types of data has led to extensive work in creating controlled-vocabularies and ontologies in multiple areas of biology. Chemical biology data poses specific challenges that require scientists from different disciplines to collaborate. Semantic wiki technology helps to merge the two strains of ontologies and scientific annotations. As a collaborative web-authoring tool, a semantic wiki makes it easier for scientists to contribute and collaborate.

PubChem is the official data repository for the MLP. The semantic wiki that is being developed as a cheminformatics tool would extract fundamental parameters of the screening data from PubChem into existing pre-defined controlled vocabularies and allow scientists to further annotate both biological, as well as chemical data. This would make high-throughput screening data more useful for gaining insight from the chemo-biological processes and would accelerate the development of new therapeutics.

## Problem Statement

Each probe project undertaken by the MLPCN goes through multiple biological assays each of which are exposed to a number of small molecules. Each biological assay can be one of three types:

1. Isolated molecular target assays
2. Cell-free multi-component assays
3. Cell or even organism based assays

The formats (e.g., plate type), controls, output, detector (e.g., imaging microscope) etc. can vary widely. In the same fashion, the sources of the compounds used in the assays and their rationale can vary widely.

To ease the task of interpretation of the results we hope to build a set of **"minimum information for chemical biology experiments"** when the results are reported out. This implies building our own set of controlled vocabularies as well as re-using existing ones. Additionally, annotating the experiments with a sufficiently dense set of controlled meta-information like ATCC number for identifying cell-lines, Uniprot Id for the protein, the GO biological process being studied will help in generating multiple ways to browse the set of experimental results.

## Problem Solution

To address these problems we are adopting a Semantic Wiki for the MLPCN scientists to use for

1. Generating appropriate controlled vocabularies
2. Annotating each experiment with suitable meta-information
3. Tools for navigating the entire set of experimental results from multiple perspectives

## Generating Appropriate Controlled Vocabularies

Semantic MediaWiki (SMW), an enhancement of the MediaWiki software, provides an open content development platform that would allows the generation of appropriate controlled vocabularies. It also provides a mechanism for collaboration between multiple stakeholders for vocabularies review, annotation, and recommendation.

The infrastructure of the SMW provides the following benefits:
- Enables property, inline queries, and new ways to explore content
- Creates a collaborative environment for collective editing and harvesting collective knowledge
- Keeps a complete record of change history
- Combines structured and unstructured knowledge for reusable contents
- Provides a formal RDF export for inferencing



Figure 1: Assays Annotation.

## Annotating Experiments with Meta-Information

**PubChem** contains thousands of chemical, biological, and bioassay screening data. The BioAssayOnDemand extension assist scientists to extract and aggregate relevant chemical biology data from PubChem for the generation of meta-data and semantic authoring and annotations. BioAssayOnDemand is a module for the PagesOnDemand extension. It generates wiki articles about bioassays in the PubChem database. Its greatest utility is in adding semantic content to the extracted data.



Figure 2: Assays extracted from PubChem for annotation.

Semantic Forms is an extension to MediaWiki that allows users to create structured forms for adding and editing data. Forms can be created on-the-fly based on existing data, the form definition and the templates that the form outputs



Figure 3: Assays Annotation with Template and Semantic Forms.

## Tools for Navigation

The Medial Subject Headings (MeSH) provides an invaluable resource to locate descriptors of possible interest and to display their hierarchical relationship. However, it does not directly connect to database retrieval systems. In order to link MeSH to other sources of biological information, a Semantic MeSH Tree Browser is implemented.

The Semantic MeSH Tree browser is similar to the MeSH Browser, except that it displays relationships between MeSH terms and assay pages in the wiki. This is possible because the semantic annotation of assays allows them to be dynamically queried via properties. For example, below we display the Semantic Tree Browser where each node displays the relationship between a specific MeSH term and Assay IDs.



Figure 4: The Semantic MeSH Tree Browser.

## References

[1] Semantic wiki: http://en.wikipedia.org/wiki/Semantic_wiki
[2] PageOnDemand: http://www.mediawiki.org/wiki/Extension:PagesOnDemand
[3] Semantic Forms: http://www.mediawiki.org/wiki/Extension:Semantic_Forms
[4] Category Tree: http://www.mediawiki.org/wiki/Extension:CategoryTree